

Table of Contents

1.0 Introduction.....	1
1.1 Background.....	1
1.2 Description of Work.....	3
2.0 ASCI Applications Overview.....	4
3.0 Vision for Tera-Scale Computing.....	7
3.1 Hardware Environment.....	7
3.2 Software Environment.....	9
4.0 SST High-Level Technical Requirements.....	11
4.1 SST Hardware High-Level Requirements.....	12
4.2 SST Software High-Level Requirements.....	19
5.0 ID High-Level Technical Requirements.....	33
5.1 Initial Delivery Hardware Requirements.....	34
5.2 Initial Delivery Software Requirements.....	37
5.3 Performance of the ID System.....	42
6.0 Implementing a Sustained Stewardship TeraFLOP.....	46
6.1 Detailed Project Plan.....	46
6.2 Project Milestones.....	50
6.3 Performance Reviews (MR).....	52
6.4 SST TeraFLOP/s sPPM Demonstration (MR).....	52
Appendix A Glossary.....	i
Hardware.....	i
Software.....	iii

Definitions

Particular paragraphs of the Statement of Work have the following designations and definitions.

- (a) **Mandatory requirements designated as (MR)**
Mandatory requirements, indicated with the verb "shall", are items that are essential to the University requirements and reflect the minimum qualifications an offeror must meet in order to have their proposal evaluated further for selection (see also Attachment 4, Evaluation Criteria).
- (b) **Mandatory Option requirements designated as (MO)**
Mandatory Option requirements deal with features, components, performance characteristics, or upgrades whose availability as an option is deemed a Mandatory Requirement by the University. Hence, a proposal not meeting a Mandatory Option will be deemed technically nonresponsive. Because the University may variously elect to include or exclude such options in resulting orders, each should appear as a separately identifiable item in the Price and Administration Proposal. Exception: Offeror need only respond to either Specification 5.1.3.1.1. or Specification 5.1.3.1.2.
- (b) **Target Requirements designated as (TR).**
Each paragraph so labeled deals with features, components, performance characteristics or other properties that is considered a desirable part of the ASCI system but will not be a determining factor of response compliance. Requirements in the Statement of Work indicated with the verb "may" are targets. Target Requirement responses will be considered as part of the evaluation of Technical Excellence (see Attachment 4, Evaluation Criteria).

1.2.10 Memory Upgrade (MO)

At the option of the University, upgrade the TeraFLOP/s scalable clusters as described in paragraph 6.2.15.

1.2.11 Performance Reviews (MR)

Apply Laboratory-provided metrics (including quarterly reviews) during the course of the contract to gauge progress.

1.2.12 Successful Project Completion (MR)

Demonstrate the achievement of the development objectives at the conclusion of the contract.

End of Section 1

the objectives of this project and Subcontractor's Research and Development Plan, which the Subcontractor believes will be of benefit to the project.

4.1 SST Hardware High-Level Requirements

4.1.1 Scalable Cluster of SMPs

4.1.1.1 Sustained Stewardship TeraFLOP SMP Scalable Cluster (MR)

The Subcontractor shall provide a Sustained Stewardship TeraFLOP/s (SST) system composed of multiple Shared memory Multi-Processors (SMPs) connected via a scalable intra-cluster communications technology. The system shall have a peak plus sustained performance of at least four (4.0) TeraFLOP/s and sustain at least one (1.0) TeraFLOP/s (1.0×10^{12} floating point operations per second) on the sPPM benchmark.

Example: If "p" is the peak performance of the system and "s" is the sustained on sPPM performance of the system, then the above mandatory requirement states that

$$s \geq 1, \text{ and } p+s \geq 4.0.$$

If we define the machine efficiency as $e = s/p$, then the above equations become:

$$p(e) \geq 31/e \quad \text{for } e < 1/3$$

and

$$p(e) \geq 34.0/(1.0+e) \quad \text{for } e \geq 1/3$$

Hence,

$$p(1/2) \geq 34/1.5 = 22.67, p(1/3) \geq 33.0 \text{ and } p(1/5) \geq 31/0.2 = 155.0$$

4.1.1.2 SST Component Scaling (MR)

In order to provide the maximum flexibility to the Subcontractor in meeting the goals of the ASCI project the exact configuration of the SST SMP scalable cluster is not specified. Rather the SST configuration is given in terms of lower bounds on component attributes. The SST SMP scalable cluster configuration shall meet or exceed the following parameters:

- Memory Size 30.5 TB
- Disk Space 375 TB
- Cache Bandwidth 312 TB/s
- Memory Bandwidth 33 TB/s
- Intra-Cluster Network Bi-Section Bandwidth 30.5 Tb/s
- System Peak Disk I/O Bandwidth 390 GB/s

4.1.1.2.1 Additional Intra-Cluster Network Bi-Section Bandwidth (TR)

The Subcontractor may provide a configuration identical to that specified in Requirement 4.1.1.2, but with an Intra-Cluster Network Bi-Section Bandwidth 31.5 Tb/s.

4.1.5.9 SMP Processor Failure Tolerance (TR)

The SMP may be able to run with one or more computational processors disabled, and to do so with minimal performance degradation. That is, the SMP may be able to tolerate failures through graceful degradation of performance.

4.1.5.10 SMP Memory Failure Tolerance (TR)

The Subcontractor may propose SMPs that are able to run with one or more memory components disabled, and to do so with minimal performance degradation. That is, the SMPs may be able to tolerate failures through graceful degradation of performance.

4.1.5.11 Replacement Parts and Maintenance (MR)

The Subcontractor shall supply hardware and software maintenance for the proposed system for a four year utilization period. Hardware maintenance response time shall be less than four hours from incident report until Subcontractor personnel arrive for repair work. Software maintenance shall include a trouble reporting mechanism and periodic software updates. In addition, the Subcontractor shall provide quick turnaround of software fixes to reported bugs. The proposed system will be installed in a classified area at the Laboratory and so maintenance personnel shall obtain DOE Q clearances.

4.2 SST Software High-Level Requirements

4.2.1 Operating System

4.2.1.1 SMP Base Operating System and License (MR)

The Subcontractor shall provide a standard multiuser POSIX (IEEE 1003.1-1990; FIPS 151-2; IEEE 1003.2 or later) compliant UNIX interactive operating system on each SMP, consisting of a basic kernel that supports system services and multiprocessing applications. Fully supported thread operations in shared address space, as defined by the POSIX 1003.1c-1995 (or later), implemented at the kernel level shall also be provided (within six months of standardization or at SST delivery). The operating system shall provide mechanisms to share memory between user processes and to run threads within a single user process on multiple CPUs simultaneously. This shall include provision of right-to-use license for an unlimited number of users, including unlimited concurrent usage, of the operating system, daemons, and associated utilities. The University will accept the Offeror's self-certification for POSIX compliance.

4.2.1.1.1 X/Open OS Compliance (TR)

The proposed operating system may have the X/Open XPG4 UNIX brand. Software with the functionality of the following X/Open components may be provided: XPG4 C Language ("ISO C"); XPG4 ISAM; FIPS 151-2; XPG4 Commands and Utilities V2; XPG4 Internationalized System Calls and Libraries (Extended); XPG4 X Window System Application Interface (FIPS 158-1). For the XPG4 UNIX brand, the

5.2.1 Operating System

5.2.1.1 SMP Base Operating System and License (MR)

The Subcontractor shall provide a standard multiuser POSIX (IEEE 1003.1-1990; FIPS 151-2; IEEE 1003.2 or later) UNIX interactive operating system on each SMP, consisting of a basic kernel that supports system services and multiprocessing applications. A fully supported thread operations in shared address space, as defined by the POSIX 1003.1c-1995 (or later), implemented at the kernel level shall also be provided. The operating system shall provide mechanisms to share memory between user processes and to run threads within a single user process on multiple CPUs simultaneously. This shall include provision of right-to-use license for an unlimited number of users, including unlimited concurrent usage, of the operating system, daemons, and associated utilities. The University will accept the Offeror's self-certification for POSIX compliance.

5.2.1.2 Networking Protocols (MR)

The operating system shall support the DoD standard networking protocol suite operating over the network interfaces described elsewhere in this document. In particular, the TCP/IP, UDP, NIS, NFS (client and server), RIP, telnet, and ftp protocols shall be supported.

5.2.1.3 Third Party Transfers (TR)

The Subcontractor may provide a driver that supports third party transfers of files between the ID and HIPPI network attached disks controlled by NSL UniTree version 2.1, or later. This implies an operating system driver capable of IPI-3 over HIPPI with National Storage Laboratory third party extensions (see LLNL Report UCRL-ID-123184).

5.2.1.4 Group Routing (MR)

The Subcontractor shall provide an implementation of "Group Routing," which segregates network traffic based on (sub)network address and group ID. A modified ROUTE table and command that allows (or explicitly disallows) packet routing to specific IP subnets based on group ID would satisfy this need.

5.2.2 Distributed Computing Environment (MR)

The Subcontractor shall provide the Open Software Foundation (OSF) Distributed Computing Environment (DCE), version 1.1 or later, client software on the proposed cluster of SMPs. This shall include Distributed File System (DFS) client-side distributed filesystem implementation which supports all standard features such as integration with CDS naming, integration with the DCE security, authentication, and authorization system. Additionally, this shall include fully supported implementation of DCE client-side security system including authentication, authorization controls, and access control lists. Fully supported remote system access programs such as rcp, rlogin, rsh, rexec, telnet, and ftp shall attempt to forward credentials to the remote system; and remote access services such

improve the hardware and software environment. Hardware technology updates may meet or exceed the following component scaling parameters:

- Memory Size/Peak FP (Byte/FLOP/s) 30.5
- Disk Space/Peak FP (Byte/FLOP/s) 325
- Cache Bandwidth/Peak FP (Byte/s/FLOP/s) 34
- Memory Bandwidth/Peak FP (Byte/s/FLOP/s) 31
- Intra-Cluster Network Bi-Section Bandwidth/Peak FP (Bits/s/FLOP/s) 30.167
- System Peak Disk I/O Bandwidth/Peak FP (Byte/s/FLOP/s) 30.03

6.2.7 FY98 Plan and Review (MR)

The Subcontractor shall provide a detailed plan of activities and deliverables for fiscal year 1998 for University review and approval in the first quarter of FY98.

6.2.8 SST Applications Development Support (MR)

The Subcontractor shall supply at least two on-site analysts to provide expertise to the University code development teams in the areas of software development tools, parallel applications libraries and applications performance at the three months prior to the SST system delivery.

6.2.9 Scalable Development Environment Demonstration (TR)

The Subcontractor may demonstrate the scalability of essential software capabilities in the application development environment across the clustered SMP system in mid CY 1998. Specifically, the Subcontractor may use the sPPM demonstration code on the full cluster to demonstrate debugger, event tracing and performance statistics capabilities.

6.2.10 Sustained Stewardship TeraFLOP (SST) Demonstration (TR)

The Subcontractor may demonstrate the SST scalable cluster in mid CY 1998 containing 0.5 Terabytes (TB) of memory and one (1.0) TeraFLOP/s of sustained performance on the sPPM demonstration code. The sum of peak and sustained performance of the SST scalable cluster shall be at least four (4.0) TeraFLOP/s.

6.2.11 Scalable Development Environment Demonstration (MR)

The Subcontractor shall demonstrate the scalability of essential software capabilities in the application development environment across the clustered SMP system no later than the end CY 1998. Specifically, the Subcontractor shall use the sPPM demonstration code on the full cluster to demonstrate debugger, event tracing and performance statistics capabilities.

6.2.12 Sustained Stewardship TeraFLOP (SST) Demonstration (MR)

The Subcontractor shall demonstrate the SST scalable cluster no later than the end CY 1998 containing 0.5 Terabytes (TB) of memory and one (1.0) TeraFLOP/s of sustained performance on the sPPM demonstration code. The sum of peak and sustained performance of the SST scalable cluster shall be at least four (4.0) TeraFLOP/s.